

# Self-Ensembling for Visual Domain Adaptation

French et al., ICLR 2018

# Outline

- This paper explores the use of self-ensembling (or teacher-student) model for domain adaptation
  - ▶ consistency regularization: minimizing the distance between student and teacher network's predictions
    - ★ student network (weights) is the current state of the model
    - ★ teacher network (weights) is the moving average of all previous states of the model
  - ▶ consistency regularization is applied on the unlabeled target data
- Ad-hoc techniques used in this model
  - ▶ confidence threshold to filter out teacher network's poor predictions
  - ▶ modulating normalization statistics of two domains

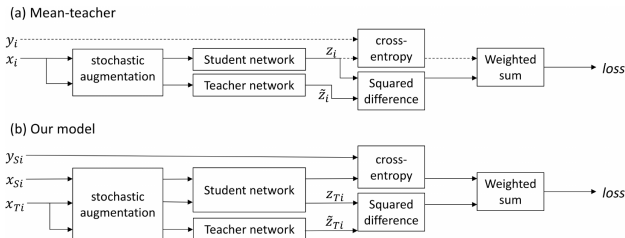
# Teacher-Student Model with Consistency Regularization

- Consistency regularization: minimizing the distance between student and teacher network's predictions

$$L_c(x_{ul}, \theta) = \mathcal{L}(p(y|x_{ul}; \theta), p^+(y|x_{ul}; \theta^+))$$

- Teacher and student networks
  - ▶ weights  $\theta$  of the student network  $p$  is the current state of the main network and is updated via back-propagation
  - ▶ weights  $\theta^+$  of the teacher network  $p^+$  is an moving average of the past states of the network
    - ★  $\theta^+ = \alpha\theta^+ + (1 - \alpha)\theta$ , where  $\theta$  is the current parameters
    - ★  $\theta^+$  is first initialized to be 0

# Teacher-Student Model for UDA



- This paper proposes a simple application of mean teacher model to UDA, without use of domain adaptation
  - ▶ minimizing the task objective (i.e. cross-entropy) with source labeled data
  - ▶ minimizing the consistency regularization with target unlabeled data
  - ▶ update student and source accordingly

# Ad-hoc Techniques for UDA

- Confidence threshold to discard poor teacher network's predictions
  - ▶ for the first few epochs, the teacher network has not accumulated enough past weights and is poor in making predictions
  - ▶ to mitigate this problem, the paper sets confidence threshold to select only confident predictions (based on the output softmax probability)
- Modulating normalization statistics of two domains
  - ▶ batch norm layers are typically used to standardize the input distribution into network layers that help improve the performance significantly
  - ▶ batch norm statistics are computed based on training dataset during training and use that statistics for normalization in inference
  - ▶ however, because of the nature of domain adaptation, the statistics of source and target domains are different from one another
  - ▶ only statistics for target domain is needed during inference, thus this paper creates separated vectors for these statistics, everything else stays the same

*Thank you !*

# Image Classification Benchmark

	USPS	MNIST	SVHN	MNIST	CIFAR	STL	Syn Digits	Syn Signs
	MNIST	USPS	MNIST	SVHN	STL	CIFAR	SVHN	GTSRB
TRAIN ON SOURCE								
SupSrc*	77.55	82.03	66.5	25.44	72.84	51.88	86.86	96.95
	$\pm 0.8$	$\pm 1.16$	$\pm 1.93$	$\pm 2.8$	$\pm 0.61$	$\pm 1.44$	$\pm 0.86$	$\pm 0.36$
SupSrc+TF	77.53	95.39	68.65	24.86	75.2	59.06	87.45	97.3
	$\pm 4.63$	$\pm 0.93$	$\pm 1.5$	$\pm 3.29$	$\pm 0.28$	$\pm 1.02$	$\pm 0.65$	$\pm 0.16$
SupSrc+TFA	91.97	96.25	71.73	28.69	75.18	59.38	87.16	98.02
	$\pm 2.15$	$\pm 0.54$	$\pm 5.73$	$\pm 1.59$	$\pm 0.76$	$\pm 0.58$	$\pm 0.85$	$\pm 0.20$
Specific aug. <sup>b</sup>	—	—	—	61.99	—	—	—	—
				$\pm 3.9$				
RevGrad <sup>a</sup> [1]	74.01	91.11	73.91	35.67	66.12	56.91	91.09	88.65
DCRN [2]	73.67	91.8	81.97	40.05	66.37	58.65	—	—
G2A [3]	90.8	92.5	84.70	36.4	—	—	—	—
ADDA [4]	90.1	89.4	76.00	—	—	—	—	—
ATT [5]	—	—	86.20	52.8	—	—	93.1	96.2
SBADA-GAN [6]	97.60	95.04	76.14	61.08	—	—	—	—
ADA [7]	—	—	97.6	—	—	—	91.86	97.66
OUR RESULTS								
MT+TF	98.07	<b>98.26</b>	99.18	13.96 <sup>c</sup>	80.08	18.3	15.94	98.63
	$\pm 2.82$	$\pm 0.11$	$\pm 0.12$	$\pm 4.41$	$\pm 0.25$	$\pm 9.03$	$\pm 0.0$	$\pm 0.09$
MT+CT*	92.35	88.14	93.33	33.87 <sup>c</sup>	77.53	71.65	96.01	98.53
	$\pm 8.61$	$\pm 0.34$	$\pm 5.88$	$\pm 4.02$	$\pm 0.11$	$\pm 0.67$	$\pm 0.08$	$\pm 0.15$
MT+CT+TF	97.28	98.13	98.64	34.15 <sup>c</sup>	79.73	<b>74.24</b>	96.51	98.66
	$\pm 2.74$	$\pm 0.17$	$\pm 0.42$	$\pm 3.56$	$\pm 0.45$	$\pm 0.46$	$\pm 0.08$	$\pm 0.12$
MT+CT+TFA	<b>99.54</b>	98.23	<b>99.26</b>	37.49 <sup>c</sup>	<b>80.09</b>	69.86	<b>97.11</b>	<b>99.37</b>
	$\pm 0.04$	$\pm 0.13$	$\pm 0.05$	$\pm 2.44$	$\pm 0.31$	$\pm 1.97$	$\pm 0.04$	$\pm 0.09$
Specific aug. <sup>b</sup>	—	—	—	<b>97.0<sup>c</sup></b>	—	—	—	—
				$\pm 0.06$				
TRAIN ON TARGET								
SupTgt*	99.53	97.29	99.59	95.7	67.75	88.86	95.62	98.49
	$\pm 0.02$	$\pm 0.2$	$\pm 0.08$	$\pm 0.13$	$\pm 2.23$	$\pm 0.38$	$\pm 0.2$	$\pm 0.32$
SupTgt+TF	99.62	97.65	99.61	96.19	70.98	89.83	96.18	98.64
	$\pm 0.04$	$\pm 0.17$	$\pm 0.04$	$\pm 0.1$	$\pm 0.79$	$\pm 0.39$	$\pm 0.09$	$\pm 0.09$
SupTgt+TFA	99.62	97.83	99.59	96.65	70.03	90.44	96.59	99.22
	$\pm 0.03$	$\pm 0.17$	$\pm 0.06$	$\pm 0.11$	$\pm 1.13$	$\pm 0.38$	$\pm 0.09$	$\pm 0.22$
Specific aug. <sup>b</sup>	—	—	—	97.16	—	—	—	—
				$\pm 0.05$				